

Technical Contribution 1: Lead Backend & ML Engineer – Contributor A

Overview

This document explains the technical contributions made by Contributor A, who served as the Lead Backend & ML Engineer for AI Shield. The contributions include developing the FastAPI backend infrastructure, implementing anomaly detection systems, building the ML pipeline, managing the Isolation Forest model, creating feature extraction modules, developing scoring systems, integrating VirusTotal and Hybrid Analysis, and implementing the complete threat detection pipeline from detection to database to WebSocket.

1. FastAPI Backend Development

What is FastAPI Backend?

FastAPI is a modern, high-performance web framework for building APIs. The backend serves as the core server that handles all requests from the frontend, processes security scans, manages the database, and provides real-time updates through WebSockets.

How It Works (Simple Explanation)

- **API Server:** Receives requests from the frontend dashboard
- **Request Processing:** Handles file uploads, scans, threat queries, and configuration
- **Service Integration:** Connects to ML models, cloud services, and security modules
- **Database Operations:** Stores and retrieves threat data, scan results, and settings
- **Real-Time Updates:** Sends live updates to connected clients via WebSocket

Key Features

- **RESTful API:** Standard HTTP endpoints for all operations
- **WebSocket Support:** Real-time bidirectional communication
- **CORS Configuration:** Secure cross-origin resource sharing
- **Modular Architecture:** Organized with routers for different features
- **Async Operations:** Non-blocking operations for better performance
- **Error Handling:** Comprehensive error handling and logging
- **File Upload Support:** Handles file uploads for scanning
- **Background Tasks:** Asynchronous job processing

Technical Implementation Points

- Uses FastAPI framework with Uvicorn ASGI server
- Modular router system (`anomaly` , `webshield` , `sandbox` , `snort` ,
`threat_actions` , `scanner` , `cloud_protection`)
- WebSocket manager for real-time connections
- SQLModel ORM for database operations
- Pydantic models for request/response validation
- CORS middleware for frontend integration
- Background task processing with asyncio

Questions & Answers

Q: Why FastAPI instead of other frameworks?

A: FastAPI provides:

- High performance (comparable to Node.js and Go)

- Automatic API documentation
- Type validation with Pydantic
- Modern Python async/await support
- Easy integration with WebSockets

Q: How does the backend handle multiple requests?

A: The backend uses:

- Async/await for non-blocking operations
- Background tasks for long-running operations
- Connection pooling for database operations
- Efficient request routing

Q: Can the backend scale to handle many users?

A: Yes. FastAPI supports:

- Multiple worker processes

- Load balancing

- Horizontal scaling

- Efficient async operations

Q: How are errors handled?

A: The backend includes:

- Try-catch blocks for exception handling

- Error logging for debugging

- User-friendly error messages

- Graceful degradation when services fail

Q: What security features are built into the backend?

A: Security features include:

- CORS configuration
 - Input validation with Pydantic
 - SQL injection prevention (SQLModel ORM)
 - File upload validation
 - Rate limiting capabilities
-

2. Anomaly Detection System

What is Anomaly Detection?

Anomaly detection identifies files or behaviors that are unusual or suspicious compared to normal patterns. Instead of looking for known malware signatures, it detects deviations from expected behavior.

How It Works (Simple Explanation)

- **Feature Extraction:** Analyzes files to extract characteristics (size, entropy, file type, etc.)
- **Pattern Analysis:** Looks for suspicious patterns (encryption, obfuscation, suspicious APIs)
- **ML Scoring:** Uses machine learning model to calculate risk score
- **Verdict Generation:** Classifies files as benign, suspicious, or malicious
- **Confidence Assessment:** Provides confidence level for the detection

Key Features

- **Multiple Detection Methods:**
 - Statistical analysis (entropy, byte patterns)
 - Heuristic rules (suspicious strings, API calls)
 - Machine learning (Isolation Forest model)
 - File structure analysis (PE headers, sections)

- **Comprehensive Feature Extraction:**

- File size and type
- Entropy calculation
- String pattern analysis
- PE/ELF structure analysis
- Import/export analysis
- Obfuscation detection

- **Risk Scoring:** Combines multiple indicators into a single risk score

- **Confidence Levels:** Provides confidence scores for detections

- **YARA Integration:** Supports YARA rule matching

Technical Implementation Points

- Uses `anomaly.py` module for all detection logic

- Isolation Forest ML model for anomaly detection

- Feature extraction from file headers and content

- Statistical analysis (entropy, byte patterns)

- Heuristic rule matching
- PE/ELF binary analysis
- String pattern detection
- Integration with YARA rules

Questions & Answers

Q: How does anomaly detection differ from signature-based detection?

A:

- **Signature-based:** Looks for known malware patterns (requires updates)
- **Anomaly-based:** Detects unusual behavior (catches new/unknown threats)
- **Combined approach:** Uses both for comprehensive protection

Q: What makes a file suspicious?

A: Files are flagged as suspicious if they have:

- High entropy (encrypted/packed)
- Suspicious API calls (process injection, registry modification)
- Obfuscated code
- Unusual file structure
- Suspicious strings (malware-related keywords)
- Abnormal file size or type

Q: Can anomaly detection have false positives?

A: Yes, but the system:

- Uses confidence scores to reduce false positives
- Combines multiple indicators for accuracy
- Allows manual review and whitelisting
- Learns from user feedback

Q: How accurate is the anomaly detection?

A: Accuracy depends on:

- Quality of training data
- Feature selection
- Model tuning
- Combination with other detection methods
- Typically achieves 85-95% accuracy with proper tuning

Q: Can the system detect zero-day threats?

A: Yes. Anomaly detection can identify:

- New malware variants
- Unknown attack patterns
- Zero-day exploits

- Polymorphic malware
-

3. ML Pipeline

What is an ML Pipeline?

An ML pipeline is a series of steps that process data through a machine learning model. It includes data collection, feature extraction, model training, prediction, and result interpretation.

How It Works (Simple Explanation)

- **Data Collection:** Gathers file samples for training
- **Feature Extraction:** Extracts relevant characteristics from files

- **Data Preprocessing:** Normalizes and scales features
- **Model Training:** Trains the Isolation Forest model on extracted features
- **Model Evaluation:** Tests model accuracy and performance
- **Model Deployment:** Saves trained model for use in production
- **Prediction:** Uses trained model to score new files

Key Features

- **Isolation Forest Algorithm:** Unsupervised learning for anomaly detection
- **Feature Engineering:** Extracts 11+ features from files
- **Data Preprocessing:** StandardScaler for feature normalization
- **Model Persistence:** Saves models as pickle files
- **Training Scripts:** Multiple training options (CSV, directory, custom)

- **Model Versioning:** Tracks model versions and configurations
- **Performance Metrics:** Tracks training accuracy and contamination rate

Technical Implementation Points

- Uses scikit-learn's IsolationForest
- Feature extraction from file metadata and content
- StandardScaler for feature normalization
- Model persistence with pickle
- Training scripts in `backend/scripts/`
- Model files stored in `backend/models/`
- Feature names stored in JSON format

Questions & Answers

Q: What is Isolation Forest?

A: Isolation Forest is:

- An unsupervised machine learning algorithm
- Designed specifically for anomaly detection
- Works by isolating anomalies in random trees
- Fast and efficient for large datasets
- Doesn't require labeled training data

Q: How is the model trained?

A: Training process:

- Collect file samples (benign and malicious)

- Extract features from each file
- Normalize features using StandardScaler
- Train Isolation Forest with contamination rate
- Save model, scaler, and feature names

Q: How often should the model be retrained?

A: Retraining should occur:

- When new threat patterns emerge
- When false positive/negative rates increase
- Periodically (monthly/quarterly)
- After collecting significant new data

Q: What features are used for training?

A: Features include:

- File size (logarithmic)
- Entropy (randomness measure)
- Non-ASCII ratio
- Printable character ratio
- File type indicators (PE, ELF, PDF, ZIP)
- Script/image detection
- Suspicious string hits

Q: Can I use my own training data?

A: Yes. The system supports:

- Training from CSV files
- Training from directory of files
- Custom feature extraction

- Configurable contamination rates
-

4. Isolation Forest Model Management

What is Isolation Forest Model Management?

Model management involves creating, training, saving, loading, and updating the Isolation Forest machine learning model used for anomaly detection.

How It Works (Simple Explanation)

- **Model Creation:** Initializes Isolation Forest with parameters
- **Training:** Trains model on file features
- **Saving:** Persists model to disk for reuse

- **Loading:** Loads saved model for predictions
- **Updating:** Retrains model with new data
- **Versioning:** Tracks different model versions

Key Features

- **Model Configuration:**

- Number of estimators (trees): 200
- Contamination rate: 0.1 (10% expected anomalies)
- Random state: 42 (for reproducibility)

- **Model Persistence:** Saves model as `.pkl` file

- **Scaler Management:** Separate scaler for feature normalization

- **Feature Tracking:** Stores feature names in JSON

- **Model Loading:** Automatic model loading on startup

- **Fallback Handling:** Graceful degradation if model unavailable

Technical Implementation Points

- Model stored in `backend/models/model.pkl`
- Scaler stored in `backend/models/scaler.pkl`
- Feature names in `backend/models/feature_names.json`
- Automatic loading in `anomaly.py`
- Training scripts in `backend/scripts/`
- Pickle serialization for model persistence

Questions & Answers

Q: What parameters are used for the Isolation Forest model?

A: Default parameters:

- `n_estimators=200` : Number of trees in the forest
- `contamination=0.1` : Expected proportion of anomalies (10%)
- `random_state=42` : For reproducible results

Q: How do I update the model?

A: Update process:

- Collect new training data
- Run training script (`train_isolation_forest.py` or `train_from_csv.py`)
- New model replaces old model files
- Restart backend to load new model

Q: What if the model file is missing?

A: The system:

- Falls back to heuristic-based detection

- Logs warning messages
- Continues operating with reduced accuracy
- Allows manual model training

Q: How large is the model file?

A: Model size depends on:

- Number of estimators (trees)
- Feature count
- Typically 1-5 MB for 200 estimators

Q: Can I use a different ML algorithm?

A: Yes, but requires:

- Modifying `anomaly.py` to use different algorithm
- Updating training scripts

- Ensuring feature compatibility

- Testing and validation

5. Feature Extraction

What is Feature Extraction?

Feature extraction is the process of analyzing files and extracting measurable characteristics (features) that can be used by machine learning models to make predictions.

How It Works (Simple Explanation)

- **File Analysis:** Reads file header and content
- **Metadata Extraction:** Gets file size, type, extension

- **Content Analysis:** Analyzes file content for patterns
- **Statistical Calculation:** Computes entropy, ratios, counts
- **Feature Vector Creation:** Combines all features into a single vector
- **Normalization:** Scales features for model compatibility

Key Features

- **11+ Features Extracted:**
 - File size (logarithmic)
 - Entropy (randomness)
 - Non-ASCII ratio
 - Printable character ratio
 - File type flags (PE, ELF, PDF, ZIP)
 - Script detection
 - Image detection
 - Suspicious string hits
- **Multiple Analysis Methods:**
 - Header analysis (magic bytes)

- Content analysis (byte patterns)
 - String analysis (text patterns)
 - Structure analysis (PE/ELF headers)
- **Efficient Processing:** Samples first 512KB for performance

- **Cross-Platform:** Works on Windows, Linux, macOS

Technical Implementation Points

- Feature extraction in `extract_features()` function
- Entropy calculation using Shannon entropy
- File type detection from magic bytes
- String pattern matching
- PE/ELF structure parsing
- Feature normalization with StandardScaler

Questions & Answers

Q: Why extract features instead of using raw file data?

A: Features provide:

- Reduced dimensionality (smaller data)
 - Meaningful characteristics
 - Normalized values
 - Better model performance
 - Faster processing
-
- Q:** How are features normalized?
- A:** Features are:
- Scaled using StandardScaler ($\text{mean}=0, \text{std}=1$)
 - Normalized to prevent feature dominance

- Applied consistently in training and prediction

Q: What if a file is too large to analyze?

A: The system:

- Samples first 512KB of file
- Analyzes file header (first 16 bytes)
- Uses file metadata (size, type)
- Provides efficient processing

Q: Can I add custom features?

A: Yes, by:

- Modifying `extract_features()` function
- Updating feature names JSON
- Retraining the model

- Ensuring feature consistency

Q: How accurate is feature extraction?

A: Feature extraction is:

- Deterministic (same file = same features)
- Fast (processes files in milliseconds)
- Reliable (handles various file types)
- Extensible (easy to add new features)

6. Scoring Modules

What are Scoring Modules?

Scoring modules calculate risk scores for files based on extracted features, heuristics, and machine learning predictions. They combine multiple indicators into a single risk assessment.

How It Works (Simple Explanation)

- **Feature Scoring:** Scores individual features
- **Heuristic Scoring:** Applies rule-based scoring
- **ML Scoring:** Uses ML model for anomaly score
- **Combination:** Combines scores from multiple sources
- **Verdict Generation:** Classifies as benign, suspicious, or malicious
- **Confidence Calculation:** Provides confidence level

Key Features

- **Multi-Source Scoring:**

- ML model score (Isolation Forest)
- Heuristic score (rule-based)
- Feature-based score (individual features)
- Pattern-based score (suspicious patterns)

- **Risk Levels:**

- **Benign:** Low risk (score < 0.3)
- **Suspicious:** Medium risk (score 0.3-0.7)
- **Malicious:** High risk (score > 0.7)

- **Confidence Scores:** Indicates reliability of detection

- **Detailed Breakdown:** Shows contribution of each indicator

Technical Implementation Points

- Main scoring function: `score_path()` in `anomaly.py`

- Combines ML prediction with heuristics

- Risk score calculation (0.0 to 1.0)

- Verdict classification
- Confidence score calculation
- Detailed result dictionary

Questions & Answers

Q: How is the final risk score calculated?

A: Final score combines:

- ML model prediction (anomaly score)
- Heuristic rules (suspicious patterns)
- Feature analysis (entropy, structure)
- Weighted combination of all factors

Q: What does a risk score of 0.8 mean?

A: A score of 0.8 indicates:

- High risk (80% suspicious)
- Likely malicious
- Should be quarantined or deleted
- Requires immediate attention

Q: Can scores be customized?

A: Yes, by:

- Adjusting ML model parameters
- Modifying heuristic rules
- Changing score thresholds
- Customizing weight combinations

Q: How are false positives handled?

A: False positives are reduced by:

- Using confidence scores
- Combining multiple indicators
- Allowing manual review
- Whitelisting safe files
- Continuous model improvement

Q: What's the difference between score and verdict?

A:

- **Score:** Numeric value (0.0-1.0) indicating risk level
- **Verdict:** Categorical classification (benign/suspicious/malicious)
- **Confidence:** Reliability of the score/verdict

7. VirusTotal Integration

What is VirusTotal?

VirusTotal is a cloud-based threat intelligence service that aggregates results from multiple antivirus engines and provides file reputation checks.

How It Works (Simple Explanation)

- **File Hash Calculation:** Computes SHA256 hash of file
- **API Request:** Sends hash to VirusTotal API
- **Reputation Check:** VirusTotal checks against known malware database
- **Result Processing:** Receives detection results from multiple engines
- **Score Calculation:** Calculates threat score based on detections

- **Integration:** Combines with local ML detection

Key Features

- **Hash-Based Lookup:** Fast reputation checks using file hashes
- **File Submission:** Automatic submission of suspicious files
- **Multi-Engine Results:** Aggregates results from 70+ antivirus engines
- **Caching:** Caches results to reduce API calls
- **Auto-Submission:** Automatically submits threats for analysis
- **Rate Limiting:** Respects API rate limits
- **Error Handling:** Graceful handling of API failures

Technical Implementation Points

- Integration in `cloud_protection.py`
- VirusTotal API v2 integration
- SHA256 hash calculation
- HTTP client with `httpx`
- Result caching (1 hour TTL)
- Auto-submission for threats
- Configuration via environment variables

Questions & Answers

Q: Do I need a VirusTotal API key?

A: Yes. VirusTotal requires:

- Free API key (limited rate: 4 requests/minute)
- Paid API key (higher rate limits)
- API key set in `VIRUSTOTAL_API_KEY` environment variable

Q: How does VirusTotal integration help?

A: VirusTotal provides:

- Access to 70+ antivirus engines
- Known malware database
- File reputation information
- Community threat intelligence
- Historical detection data

Q: What happens if VirusTotal API is unavailable?

A: The system:

- Falls back to local ML detection

- Logs API errors

- Continues operating normally

- Retries on next scan

Q: Are files automatically submitted to VirusTotal?

A: Yes, if:

- Auto-submit is enabled

- File is detected as suspicious/malicious

- File size is within limits (32MB default)

- API key is configured

Q: How long does a VirusTotal check take?

A: Check duration:

- Hash lookup: 1-2 seconds
 - File submission: 5-30 seconds
 - Results depend on API response time
-

8. Hybrid Analysis Integration

What is Hybrid Analysis?

Hybrid Analysis is a cloud-based malware analysis platform that provides automated sandbox analysis and behavioral detection of suspicious files.

How It Works (Simple Explanation)

- **File Submission:** Uploads suspicious files to Hybrid Analysis

- **Sandbox Execution:** Hybrid Analysis runs file in isolated environment
- **Behavior Analysis:** Monitors file behavior (API calls, network activity)
- **Report Generation:** Creates detailed analysis report
- **Result Integration:** Integrates results with local detection

Key Features

- **Automated Sandbox:** Runs files in isolated Windows environment
- **Behavioral Analysis:** Tracks API calls and system interactions
- **Threat Intelligence:** Provides threat classification
- **File Submission:** Automatic submission of threats
- **Job Tracking:** Tracks submission status and results

- **API Integration:** Uses Hybrid Analysis API v2
- **Error Handling:** Graceful handling of API failures

Technical Implementation Points

- Integration in `cloud_protection.py`
- Hybrid Analysis API v2
- File upload with multipart form data
- Job ID tracking
- SHA256 hash calculation
- Auto-submission for threats
- Configuration via environment variables

Questions & Answers

Q: Do I need Hybrid Analysis API credentials?

A: Yes. Hybrid Analysis requires:

- API key (`HYBRID_ANALYSIS_API_KEY`)
- API secret (`HYBRID_ANALYSIS_API_SECRET`)
- Free account available with limitations

Q: What does Hybrid Analysis provide that local detection doesn't?

A: Hybrid Analysis provides:

- Real execution environment
- Behavioral analysis
- Network activity monitoring
- Detailed threat reports

- Community threat intelligence

Q: How long does Hybrid Analysis take?

A: Analysis duration:

- Submission: 5-10 seconds
- Analysis: 1-5 minutes (depending on file)
- Results available via API

Q: Are files automatically submitted?

A: Yes, if:

- Auto-submit is enabled
- File is detected as suspicious/malicious
- File size is within limits
- API credentials are configured

Q: What information does Hybrid Analysis provide?

A: Hybrid Analysis provides:

- Threat verdict (malicious/benign/suspicious)
- Behavioral indicators
- API calls made
- Network connections
- File modifications
- Registry changes

9. Threat Pipeline (Detection → DB → WebSocket)

What is the Threat Pipeline?

The threat pipeline is the complete flow from detecting a threat to storing it in the database to notifying connected clients in real-time via WebSocket.

How It Works (Simple Explanation)

- **Detection:** File is scanned and threat is detected
- **Threat Creation:** Threat record is created with details
- **Database Storage:** Threat is saved to SQLite database
- **WebSocket Broadcast:** Threat is broadcast to all connected clients
- **Frontend Update:** Dashboard receives and displays threat
- **User Action:** User can quarantine, delete, or allow threat

Key Features

- **End-to-End Flow:** Complete pipeline from detection to notification
- **Real-Time Updates:** Instant notifications via WebSocket
- **Database Persistence:** All threats stored for history
- **Asynchronous Processing:** Non-blocking operations
- **Error Handling:** Graceful error handling at each step
- **Status Tracking:** Tracks threat status (detected, quarantined, deleted, allowed)
- **Multi-Source Support:** Handles threats from ML, Snort, WebShield, Sandbox

Technical Implementation Points

- Threat detection in `anomaly.py`

- Database operations in `store.py`

- WebSocket manager in `main.py`

- Async task processing

- Threat model (SQLModel)

- Real-time broadcasting

- Frontend WebSocket client

Questions & Answers

Q: How fast is the threat pipeline?

A: Pipeline speed:

- Detection: 10-100ms (depending on file size)

- Database storage: <10ms

- WebSocket broadcast: <5ms
- Total: Typically <200ms end-to-end

Q: What happens if the database is unavailable?

A: The system:

- Logs errors
- Continues detection
- Queues threats for later storage
- Notifies administrators

Q: How are WebSocket connections managed?

A: WebSocket management:

- Connection pool in WSManager
- Automatic reconnection handling

- Dead connection cleanup
- Broadcast to all connected clients

Q: Can threats be filtered or searched?

A: Yes. The system supports:

- Filtering by severity
- Filtering by source
- Filtering by action status
- Search by file path or URL
- Pagination for large result sets

Q: What information is stored for each threat?

A: Threat records include:

- Unique ID

- Timestamp
- Severity (low/medium/high/critical)
- Description
- Source (ML/Snort/WebShield/Sandbox)
- Action status
- File path or URL
- Deep analysis data (optional)

Summary

Contributor A's technical contributions form the foundation of AI Shield's detection and analysis capabilities:

- **FastAPI Backend:** Core server infrastructure with REST API and WebSocket support

- **Anomaly Detection:** Multi-layered detection system combining ML and heuristics
- **ML Pipeline:** Complete machine learning workflow from training to prediction
- **Isolation Forest Model:** Unsupervised anomaly detection model
- **Feature Extraction:** Comprehensive feature extraction from files
- **Scoring Modules:** Risk scoring and verdict generation
- **VirusTotal Integration:** Cloud-based threat intelligence
- **Hybrid Analysis Integration:** Automated sandbox analysis
- **Threat Pipeline:** End-to-end flow from detection to real-time notification

All these systems work together to provide comprehensive threat detection, with machine learning at the core, enhanced by cloud intelligence, and delivered through a modern, scalable backend architecture.

Technical Stack

- **Backend Framework:** FastAPI 0.115+
- **ASGI Server:** Uvicorn
- **Database:** SQLite with SQLModel ORM
- **Machine Learning:** scikit-learn (Isolation Forest)
- **HTTP Client:** httpx
- **WebSocket:** FastAPI WebSocket support
- **Data Processing:** NumPy, pickle
- **File Analysis:** LIEF (for binary analysis)
- **Platform:** Python 3.8+

Document generated for AI Shield Project - Technical Contribution Documentation